

## RESEARCH

## Open Access

# Perturbation analysis of the stochastic algebraic Riccati equation

Chun-Yueh Chiang<sup>1</sup>, Hung-Yuan Fan<sup>2</sup>, Matthew M Lin<sup>3\*</sup> and Hsin-An Chen<sup>3</sup>\*Correspondence:  
mhlin@ccu.edu.tw<sup>3</sup>Department of Mathematics,  
National Chung Cheng University,  
Chia-Yi 621, Taiwan  
Full list of author information is  
available at the end of the article

## Abstract

In this paper we study a general class of stochastic algebraic Riccati equations (SARE) arising from the indefinite linear quadratic control and stochastic  $H_\infty$  problems. Using the Brouwer fixed point theorem, we provide sufficient conditions for the existence of a stabilizing solution of the perturbed SARE. We obtain a theoretical perturbation bound for measuring accurately the relative error in the exact solution of the SARE. Moreover, we slightly modify the condition theory developed by Rice and provide explicit expressions of the condition number with respect to the stabilizing solution of the SARE. A numerical example is applied to illustrate the sharpness of the perturbation bound and its correspondence with the condition number.

**MSC:** Primary 15A24; 65F35; secondary 47H10; 47H14

**Keywords:** Brouwer fixed-point theorem; perturbation bound; stochastic algebraic Riccati equations; condition number

## 1 Introduction

In this paper we consider a general class of continuous-time stochastic algebraic Riccati equations

$$A^\top X + XA + C^\top XC - (XB + C^\top XD + S)(R + D^\top XD)^{-1}(B^\top X + D^\top XC + S^\top) + H = 0, \quad (1a)$$

$$R + D^\top XD \succ 0, \quad (1b)$$

where  $A \in \mathbb{R}^{n \times n}$ ,  $C \in \mathbb{R}^{n \times n}$ ,  $B \in \mathbb{R}^{n \times m}$ ,  $D \in \mathbb{R}^{n \times m}$ ,  $S \in \mathbb{R}^{n \times m}$ , respectively. Moreover,  $H \in \mathbb{R}^{n \times n}$  and  $R \in \mathbb{R}^{m \times m}$  are symmetric matrices. Here we denote  $M \succ 0$  (respectively,  $M \succeq 0$ ) if  $M$  is symmetric positive definite (respectively, positive semidefinite). The unknown  $X \in \mathbb{R}^{n \times n}$  is a symmetric solution to SARE (1a)-(1b). Let  $\mathcal{S}^n$  be the set of all symmetric  $n \times n$  real matrices. For any  $X, Y \in \mathcal{S}^n$ , we write  $X \succeq Y$  if  $X - Y \succeq 0$ .

In essence, SARE (1a)-(1b) is a rational Riccati-type matrix equation associated with the operator  $\mathcal{R} : \text{dom } \mathcal{R} \rightarrow \mathcal{S}^n$

$$\mathcal{R}(X) = \mathcal{P}(X) - \mathcal{S}(X)\mathcal{Q}(X)^{-1}\mathcal{S}(X)^\top,$$

where the affine linear operators  $\mathcal{P} : \mathcal{S}^n \rightarrow \mathcal{S}^n$ ,  $\mathcal{Q} : \mathcal{S}^n \rightarrow \mathcal{S}^m$ ,  $\mathcal{S} : \mathcal{S}^n \rightarrow \mathbb{R}^{n \times m}$ , and  $\text{dom } \mathcal{R}$  are defined by

$$\mathcal{P}(X) = A^\top X + XA + C^\top XC + H,$$

$$\begin{aligned}\mathcal{Q}(X) &= R + D^\top X D, \\ \mathcal{S}(X) &= X B + C^\top X D + S, \\ \text{dom } \mathcal{R} &= \{X \in \mathcal{S}^n \mid \mathcal{Q}(X) \succ 0\}.\end{aligned}$$

We say that  $X$  is the maximal solution (or the greatest solution) of SARE (1a)-(1b) if it satisfies (1a)-(1b) and  $X \succeq P$  for any  $P \in \mathcal{S}^n$  satisfying  $\mathcal{R}(P) \geq 0$  and (1b), *i.e.*,  $X$  is the maximal solution of  $\mathcal{R}(X) \geq 0$  with the constraint (1b). Furthermore, it is easily seen that SARE (1a)-(1b) also contains the continuous-time algebraic Riccati equation (CARE)

$$A^\top X + X A - X B R^{-1} B^\top X + H = 0 \quad (2)$$

with  $R \succ 0$ ,  $C = 0$ ,  $D = 0$  and  $S = 0$ , and the discrete-time algebraic Riccati equation (DARE)

$$X - C^\top X C + (C^\top X D + S)(R + D^\top X D)^{-1}(D^\top X C + S^\top) - H = 0 \quad (3)$$

with  $A = \frac{A}{2}$  and  $B = 0$ , as special cases.

Matrix equations of the type (1a)-(1b) are encountered in the *indefinite* linear quadratic (LQ) control problem [1], and the disturbance attenuation problem, which is in deterministic case the  $H_\infty$  control theory, for linear stochastic systems with both state- and input-dependent white noise. For example, see [2-4]. For simplicity, we only consider one-dimensional Wiener process of white noise in this paper; it is straightforward but tedious to extend all perturbation results presented in this paper for multi-dimensional cases. In the aforementioned applications of linear stochastic systems, a symmetric solution  $X$ , called a *stabilizing solution*, to SARE (1a)-(1b) ought to be determined for the design of optimal controllers. This stabilizing solution plays a very important role in many applications of linear system control theory. The definition of a stabilizing solution to SARE (1a)-(1b) is given as follows. (See also [3, Definition 5.2].)

**Definition 1.1** Let  $X \in \mathcal{S}^n$  be a solution to SARE (1a)-(1b),  $\Phi = A + B F$  and  $\Psi = C + D F$ , where  $F = -\mathcal{Q}(X)^{-1} \mathcal{S}(X)^\top$ . The matrix  $X$  is called a stabilizing solution for  $\mathcal{R}$  if the spectrum of the associated operator  $\mathcal{L}_c$  with respect to  $X$  defined by

$$\mathcal{L}_c(W) = \Phi^\top W + W \Phi + \Psi^\top W \Psi, \quad W \in \mathcal{S}^n, \quad (4)$$

is contained in the open left half plane, *i.e.*,  $\sigma(\mathcal{L}_c) \subset \mathbb{C}_-$ .

Note that if  $C = D = 0$  in (1a)-(1b), then it is easily seen from Definition 1.1 that the matrix  $X \in \mathcal{S}^n$  is a stabilizing solution to SARE (1a)-(1b) or, equivalently, CARE (2) if and only if  $\sigma(\Phi) \subset \mathbb{C}_-$ . Therefore, Definition 1.1 is a natural generalization of the definition of a stabilizing solution to CARE (2) in classical linear control theory. Moreover, a necessary and sufficient condition for the existence of the stabilizing solution to a more general SARE is derived in Theorem 7.2 of [3]. See also [1, Theorem 10]. In this case, it is also shown that if SARE (1a)-(1b) has a stabilizing solution  $X \in \text{dom}(\mathcal{R})$ , then it is necessarily a maximal solution and thus unique [1, 3].

The standard CARE (2) and DARE (3) are widely studied and play very important roles in both classical LQ and  $H_\infty$  control problems for deterministic linear systems [5–7]. In the past four decades, an extensive amount of numerical methods were studied and developed for solving the CARE and DARE (see [8–10] and the references therein). There are two major methodologies among these numerical methods or algorithms. One is the so-called *Schur method* or *invariant subspace method*, which was first proposed by Laub [11]. According to this methodology, the unique and non-negative definite stabilizing solution of the CARE (or DARE) can be obtained by computing the stable invariant subspace (or deflating subspace) of the associated Hamiltonian matrix (or symplectic matrix pencil). Some variants of the invariant subspace method, which preserve the structure of the Hamiltonian matrix (or symplectic matrix pencil) by special orthogonal transformations in the whole computational process, are considered by Mehrmann and his coauthors [12–18]. The other methodology comes from the *iterative method*, for example, it is referred to as Newton's method [6], matrix sign function method [19], disk function method [20], and structured doubling algorithms [21, 22] and references therein. So far there has been no sources in applying the invariant subspace methods for solving SARE (1a)-(1b), since the structures of associated Hamiltonian matrix or symplectic matrix pencil are not available. Only the iterative methods, *e.g.*, Newton's method [3] and the interior-point algorithm presented in [1], can be applied to computing the numerical solutions of SARE (1a)-(1b). Recently, normwise residual bounds were proposed for assessing the accuracy of a computed solution to SARE (1a)-(1b) [23].

Due to the effect of roundoff errors or the measurement errors of experimental data, small perturbations are often incorporated in the coefficient matrices of SARE (1a)-(1b), and hence we obtain the perturbed SARE

$$\begin{aligned} \tilde{A}^\top \tilde{X} + \tilde{X} \tilde{A} + \tilde{C}^\top \tilde{X} \tilde{C} - (\tilde{X} \tilde{B} + \tilde{C}^\top \tilde{X} \tilde{D} + \tilde{S})(\tilde{R} + \tilde{D}^\top \tilde{X} \tilde{D})^{-1}(\tilde{B}^\top \tilde{X} + \tilde{D}^\top \tilde{X} \tilde{C} + \tilde{S}^\top) \\ + \tilde{H} = 0, \end{aligned} \quad (5a)$$

$$\tilde{R} + \tilde{D}^\top \tilde{X} \tilde{D} \succ 0, \quad (5b)$$

where  $\tilde{A}$ ,  $\tilde{B}$ ,  $\tilde{C}$ ,  $\tilde{D}$ ,  $\tilde{H}$ ,  $\tilde{R}$  and  $\tilde{S}$  are perturbed coefficient matrices of compatible sizes. The main question is under what conditions perturbed SARE (5a)-(5b) still has a stabilizing solution  $\tilde{X} \in \mathcal{S}^n$ . Moreover, how sensitive is the stabilizing solution  $X \in \text{dom}(\mathcal{R})$  of original SARE (1a)-(1b) with respect to small changes in the coefficient matrices? This is related to the conditioning of SARE (1a)-(1b). Therefore, we will try to answer these questions for SARE (1a)-(1b) in this paper. For CARE (2) and DARE (3), the normwise non-local and local perturbation bounds have been widely studied in the literature. See, *e.g.*, [24–26]. Also, computable residual bounds were derived for measuring the accuracy of a computed solution to CARE (2) and DARE (3), respectively [27, 28]. To our best knowledge, these issues have not been taken into account for constrained SARE (1a)-(1b) in the literature.

To facilitate our discussion, we use  $\|\cdot\|_F$  to denote the Frobenius norm and  $\|\cdot\|$  to denote the operator norm induced by the Frobenius norm. For  $A = (A_1, \dots, A_n) = (a_{ij}) \in \mathbb{R}^{m \times n}$  and  $B \in \mathbb{R}^{p \times q}$ , the Kronecker product of  $A$  and  $B$  is defined by  $A \otimes B = (a_{ij}B) \in \mathbb{R}^{mp \times nq}$ , and the operator  $\text{vec}(A)$  is denoted by  $\text{vec}(A) = (A_1^\top, \dots, A_n^\top)^\top$ . It is known that

$$\text{vec}(ABC) = (C^\top \otimes A) \text{vec}(B), \quad \text{vec}(A^\top) = P_{n,m} \text{vec}(A),$$

where  $A \in \mathbb{R}^{n \times m}$ ,  $B \in \mathbb{R}^{m \times \ell}$ ,  $C \in \mathbb{R}^{\ell \times k}$ , and  $P_{n,m}$  is the Kronecker permutation matrix which maps  $\text{vec}(A)$  into  $\text{vec}(A^\top)$  for a rectangle matrix  $A$ , i.e.,

$$P_{n,m} = \sum_{i,j=1}^{n,m} E_{i,j,n \times m} \otimes E_{j,i,m \times n},$$

where the  $n \times m$  matrix  $E_{i,j,n \times m}$  has 1 as its  $(i,j)$  entry and 0's elsewhere.

This paper is organized as follows. In Section 2, a perturbation equation is derived from SAREs (1a)-(1b) and (5a)-(5b) without dropping any higher-order terms. By using Brouwer fixed point theorem, we obtain a perturbation bound for the stabilizing solution of SARE (5a)-(5b) in Section 3. In order to guarantee the existence of the stabilizing solution of perturbed SARE (5a)-(5b), some stability analysis of the operator  $\mathcal{L}_c$  is established in Section 4. A theoretical formula of the normwise condition number of the stabilizing solution to SARE (1a)-(1b) is derived in Section 5. Finally, in Section 6, a numerical example is given to illustrate the sharpness and tightness of our perturbation bounds, and Section 7 concludes the paper.

## 2 Perturbation equation

Assume that  $X \in \mathcal{S}^n$  is the unique stabilizing solution to SARE (1a)-(1b) and  $\tilde{X} \in \mathcal{S}^n$  is a symmetric solution of perturbed SARE (5a)-(5b), that is,

$$\mathcal{R}(X) := A^\top X + XA + C^\top XC - \Xi(X) + H = 0, \quad (6)$$

$$\tilde{\mathcal{R}}(\tilde{X}) := \tilde{A}^\top \tilde{X} + \tilde{X}\tilde{A} + \tilde{C}^\top \tilde{X}\tilde{C} - \tilde{\Xi}(\tilde{X}) + \tilde{H} = 0, \quad (7)$$

where the two operator  $\Xi : \mathcal{S}^n \rightarrow \mathcal{S}^n$  and  $\tilde{\Xi} : \mathcal{S}^n \rightarrow \mathcal{S}^n$  are given by

$$\begin{aligned} \Xi(X) &= \mathcal{S}(X)\mathcal{Q}(X)^{-1}\mathcal{S}(X)^\top, \\ \tilde{\Xi}(\tilde{X}) &= \tilde{\mathcal{S}}(\tilde{X})\tilde{\mathcal{Q}}(\tilde{X})^{-1}\tilde{\mathcal{S}}(\tilde{X})^\top, \end{aligned} \quad (8)$$

and two affine linear operators  $\tilde{\mathcal{S}} : \mathcal{S}^n \rightarrow \mathcal{S}^n$ ,  $\tilde{\mathcal{Q}} : \mathcal{S}^n \rightarrow \mathcal{S}^m$  are defined by

$$\begin{aligned} \tilde{\mathcal{S}}(\tilde{X}) &= \tilde{X}\tilde{B} + \tilde{C}^\top \tilde{X}\tilde{D} + \tilde{S}, \\ \tilde{\mathcal{Q}}(\tilde{X}) &= \tilde{R} + \tilde{D}^\top \tilde{X}\tilde{D} \end{aligned}$$

for all  $\tilde{X} \in \mathcal{S}^n$ . Let

$$\Delta X = \tilde{X} - X.$$

The purpose of this section is to derive a perturbation equation of  $\Delta X$  from SAREs (1a)-(1b) and (5a)-(5b). For the sake of perturbation analysis, we adopt the following notations:

$$\begin{aligned} \Delta A &= \tilde{A} - A, & \Delta B &= \tilde{B} - B, & \Delta C &= \tilde{C} - C, & \Delta D &= \tilde{D} - D, \\ \Delta S &= \tilde{S} - S, & \Delta R &= \tilde{R} - R, & \Delta H &= \tilde{H} - H \end{aligned} \quad (9)$$

and

$$\begin{aligned} \delta Q &= \Delta R + D^\top X \Delta D + \Delta D^\top X D + \Delta D^\top X \Delta D, \\ \delta S &= \Delta S + C^\top X \Delta D + \Delta C^\top X D + \Delta C^\top X \Delta D + X \Delta B. \end{aligned} \quad (10)$$

Moreover, let

$$\begin{aligned} F &= -Q(X)^{-1}S(X)^{\top}, & \tilde{F} &= -\tilde{Q}(X)^{-1}\tilde{S}(X)^{\top}, \\ \Psi &= C + DF, & \tilde{\Psi} &= \tilde{C} + \tilde{D}\tilde{F}, \\ \Phi &= A + BF, & \tilde{\Phi} &= \tilde{A} + \tilde{B}\tilde{F}, \end{aligned} \quad (11)$$

and by the definition of  $\Psi$ , we define

$$K := X\Psi. \quad (12)$$

Note that  $\tilde{S}(X) = S(X) + \delta S$  and  $\tilde{Q}(X) = Q(X) + \delta Q$ . Substituting (11) into (8), we observe that

$$\begin{aligned} \Xi(X) &= -S(X)F, \\ \tilde{\Xi}(\tilde{X}) &= (\tilde{S}(X) + \Delta X\tilde{B} + \tilde{C}^{\top}\Delta X\tilde{D})(\tilde{Q}(X) + \tilde{D}^{\top}\Delta X\tilde{D})^{-1} \\ &\quad \times (\tilde{S}(X) + \Delta X\tilde{B} + \tilde{C}^{\top}\Delta X\tilde{D})^{\top}. \end{aligned} \quad (13)$$

Thus far, we have not specified the relation between  $\mathcal{R}(X)$  and  $\tilde{\mathcal{R}}(\tilde{X})$ . Such a tedious task can be turned into a breeze by repeatedly applying the matrix identities [29]

$$(I + U)^{-1} = I - U(I + U)^{-1}, \quad V(I + UV)^{-1} = (I + VU)^{-1}V. \quad (14)$$

To begin with, assume that  $\Delta R$  and  $\Delta D$  are sufficiently small so that  $\tilde{Q}(X)$  is invertible. We see that the product

$$\begin{aligned} &(\tilde{S}(X) + \Delta X\tilde{B} + \tilde{C}^{\top}\Delta X\tilde{D})(\tilde{Q}(X) + \tilde{D}^{\top}\Delta X\tilde{D})^{-1} \\ &= (\tilde{S}(X) + \Delta X\tilde{B} + \tilde{C}^{\top}\Delta X\tilde{D}) \\ &\quad \times [I - \tilde{Q}(X)^{-1}\tilde{D}^{\top}\Delta X\tilde{D}(I + \tilde{Q}(X)^{-1}\tilde{D}^{\top}\Delta X\tilde{D})^{-1}]\tilde{Q}(X)^{-1} \\ &= -\tilde{F}^{\top} + \Delta X\tilde{B}(I + \tilde{Q}(X)^{-1}\tilde{D}^{\top}\Delta X\tilde{D})^{-1}\tilde{Q}(X)^{-1} \\ &\quad + \tilde{\Psi}^{\top}\Delta X\tilde{D}(I + \tilde{Q}(X)^{-1}\tilde{D}^{\top}\Delta X\tilde{D})^{-1}\tilde{Q}(X)^{-1}. \end{aligned}$$

It follows that

$$\begin{aligned} \tilde{\Xi}(\tilde{X}) &= -\tilde{S}(X)\tilde{F} - \tilde{F}^{\top}\tilde{B}^{\top}\Delta X - \tilde{F}^{\top}\tilde{D}^{\top}\Delta X\tilde{C} - \tilde{\Psi}^{\top}\Delta X\tilde{D}\tilde{F} - \Delta X\tilde{B}\tilde{F} \\ &\quad + (\tilde{\Psi}^{\top}\Delta X\tilde{D} + \Delta X\tilde{B})(I + \tilde{Q}(X)^{-1}\tilde{D}^{\top}\Delta X\tilde{D})^{-1}\tilde{Q}(X)^{-1} \\ &\quad \times (\tilde{D}^{\top}\Delta X\tilde{\Psi} + \tilde{B}^{\top}\Delta X) \end{aligned}$$

since  $\tilde{F}^{\top}\tilde{S}(X)^{\top} = \tilde{S}(X)\tilde{F}$ . Next, from (11) we can see that

$$\begin{aligned} \tilde{\Phi}^{\top}\Delta X + \Delta X\tilde{\Phi} &= \tilde{A}^{\top}\Delta X + \Delta X\tilde{A} + \tilde{F}^{\top}\tilde{B}^{\top}\Delta X + \Delta X\tilde{B}\tilde{F}, \\ \tilde{\Psi}^{\top}\Delta X\tilde{\Psi} &= \tilde{C}^{\top}\Delta X\tilde{C} + \tilde{F}^{\top}\tilde{D}^{\top}\Delta X\tilde{C} + \tilde{\Psi}^{\top}\Delta X\tilde{D}\tilde{F}. \end{aligned} \quad (15)$$

Applying (15), we obtain the linear equation

$$\tilde{\mathcal{R}}(\tilde{X}) - \mathcal{R}(X) = \tilde{\Phi}^\top \Delta X + \Delta X \tilde{\Phi} + \tilde{\Psi}^\top \Delta X \tilde{\Psi} - E - h_2(\Delta X) = 0, \quad (16)$$

where

$$\begin{aligned} E &:= -(\Delta A^\top X + X \Delta A + \tilde{C}^\top X \tilde{C} - C^\top X C + \tilde{S}(X) \tilde{F} - S(X) F + \Delta H), \\ h_2(\Delta X) &:= \tilde{\Psi}^\top \Delta X \tilde{D} (I + \tilde{Q}(X)^{-1} \tilde{D}^\top \Delta X \tilde{D})^{-1} \tilde{Q}(X)^{-1} \tilde{D}^\top \Delta X \tilde{\Psi} \\ &\quad + \Delta X \tilde{B} (I + \tilde{Q}(X)^{-1} \tilde{D}^\top \Delta X \tilde{D})^{-1} \tilde{Q}(X)^{-1} \tilde{B}^\top \Delta X \\ &\quad + \Delta X \tilde{B} (I + \tilde{Q}(X)^{-1} \tilde{D}^\top \Delta X \tilde{D})^{-1} \tilde{Q}(X)^{-1} \tilde{D}^\top \Delta X \tilde{\Psi} \\ &\quad + \tilde{\Psi}^\top \Delta X \tilde{D} (I + \tilde{Q}(X)^{-1} \tilde{D}^\top \Delta X \tilde{D})^{-1} \tilde{Q}(X)^{-1} \tilde{B}^\top \Delta X. \end{aligned}$$

It follows from (16) that

$$\tilde{\Phi}^\top \Delta X + \Delta X \tilde{\Phi} + \tilde{\Psi}^\top \Delta X \tilde{\Psi} = E + h_2(\Delta X). \quad (17)$$

Equipped with this fact, we now are going to derive a perturbation equation in terms of  $\Delta X$  by using  $\Delta A$ ,  $\Delta B$ ,  $\Delta C$ ,  $\Delta D$ ,  $\Delta S$ ,  $\Delta R$ ,  $\delta S$ , and  $\delta Q$ . It should be noted that

$$\begin{aligned} \tilde{\Psi} &= (\Delta C + C) - (\Delta D + D)(Q(X) + \delta Q)^{-1}(S(X) + \delta S)^\top \\ &= \Psi + \Delta \Psi, \end{aligned}$$

with

$$\begin{aligned} \Delta \Psi &:= \Delta C - \Delta D Q(X)^{-1} S(X)^\top - \Delta D Q(X)^{-1} \delta S^\top - D Q(X)^{-1} \delta S^\top \\ &\quad + (\Delta D + D) Q(X)^{-1} \delta Q Q(X)^{-1} (I + Q(X)^{-1} \delta Q)^{-1} (S(X)^\top + \delta S^\top) \end{aligned} \quad (18)$$

and

$$\begin{aligned} \tilde{\Phi} &= (\Delta A + A) - (\Delta B + B)(Q(X) + \delta Q)^{-1}(S(X) + \delta S)^\top \\ &= \Phi + \Delta \Phi, \end{aligned}$$

with

$$\begin{aligned} \Delta \Phi &:= \Delta A - \Delta B Q(X)^{-1} S(X)^\top - \Delta B Q(X)^{-1} \delta S^\top - B Q(X)^{-1} \delta S^\top \\ &\quad + (\Delta B + B) Q(X)^{-1} \delta Q Q(X)^{-1} (I + Q(X)^{-1} \delta Q)^{-1} (S(X)^\top + \delta S^\top). \end{aligned} \quad (19)$$

It then is natural to express the left-hand side of (17) by  $\Delta \Phi$  and  $\Delta \Psi$  such that

$$\tilde{\Phi}^\top \Delta X + \Delta X \tilde{\Phi} + \tilde{\Psi}^\top \Delta X \tilde{\Psi} = \Phi^\top \Delta X + \Delta X \Phi + \Psi^\top \Delta X \Psi - h_1(\Delta X),$$

with

$$h_1(\Delta X) := -(\Delta \Phi^\top \Delta X + \Delta X \Delta \Phi + \Psi^\top \Delta X \Delta \Psi + \Delta \Psi^\top \Delta X \Psi + \Delta \Psi^\top \Delta X \Delta \Psi).$$

Observe further that

$$\begin{aligned}\tilde{C}^\top X \tilde{C} - C^\top X C &= C^\top X \Delta C + \Delta C^\top X C + \Delta C^\top X \Delta C, \\ \tilde{S}(X) \tilde{F} - S(X) F &= -(S(X) + \delta S)(Q(X) + \delta Q)^{-1}(S(X) + \delta S)^\top \\ &\quad + S(X) Q(X)^{-1} S(X)^\top \\ &= F^\top \delta S^\top + \delta S F - \delta S(I + Q(X)^{-1} \delta Q)^{-1} Q(X)^{-1} \delta S^\top \\ &\quad - F^\top \delta Q(I + Q(X)^{-1} \delta Q)^{-1} Q(X)^{-1} \delta S^\top \\ &\quad - \delta S Q(X)^{-1} \delta Q(I + Q(X)^{-1} \delta Q)^{-1} F \\ &\quad + F^\top \delta Q F - F^\top \delta Q(I + Q(X)^{-1} \delta Q)^{-1} Q(X)^{-1} \delta Q F.\end{aligned}$$

Upon substituting (10) into  $\delta S F$  and  $F^\top \delta Q F$ , we have

$$\begin{aligned}\delta S F &= \Delta S F + C^\top X \Delta D F + \Delta C^\top X D F + \Delta C^\top X \Delta D F + X \Delta B F, \\ F^\top \delta Q F &= F^\top \Delta R F + F^\top D^\top X \Delta D F + F^\top \Delta D^\top X D F + F^\top \Delta D^\top X \Delta D F,\end{aligned}$$

so that the structure of  $E$  in (17) can be partitioned into linear equations

$$\begin{aligned}E_1 &:= -(K^\top \Delta D F + F^\top \Delta D^\top K + K^\top \Delta C + \Delta C^\top K + F^\top \Delta R F + F^\top \Delta S^\top + \Delta S F + \Delta H), \\ E_2 &:= -[\Delta A^\top X + X \Delta A + \Delta C^\top X \Delta C + F^\top \Delta B^\top X + X \Delta B F \\ &\quad + F^\top \Delta D^\top X \Delta C + \Delta C^\top X \Delta D F + F^\top \Delta D^\top X \Delta D F \\ &\quad - (F^\top \delta Q + \delta S)(I + Q(X)^{-1} \delta R)^{-1} Q(X)^{-1} (\delta Q F + \delta S^\top)],\end{aligned}$$

that is,  $E = E_1 + E_2$ .

**Lemma 2.1** *Let  $X$  be the stabilizing solution of SARE (1a)-(1b) and  $\tilde{X}$  be a symmetric solution of perturbed SARE (5a)-(5b). If  $\Delta X = \tilde{X} - X$ , then  $\Delta X$  satisfies the equation*

$$\Phi^\top \Delta X + \Delta X \Phi + \Psi^\top \Delta X \Psi = E_1 + E_2 + h_1(\Delta X) + h_2(\Delta X), \quad (20)$$

where

$$\begin{aligned}E_1 &= -(K^\top \Delta D F + F^\top \Delta D^\top K + K^\top \Delta C + \Delta C^\top K + F^\top \Delta R F \\ &\quad + F^\top \Delta S^\top + \Delta S F + \Delta H),\end{aligned} \quad (21a)$$

$$\begin{aligned}E_2 &= -[\Delta A^\top X + X \Delta A + \Delta C^\top X \Delta C + F^\top \Delta B^\top X + X \Delta B F \\ &\quad + F^\top \Delta D^\top X \Delta C + \Delta C^\top X \Delta D F + F^\top \Delta D^\top X \Delta D F \\ &\quad - (F^\top \delta Q + \delta S)(I + Q(X)^{-1} \delta R)^{-1} Q(X)^{-1} (\delta Q F + \delta S^\top)],\end{aligned}$$

$$h_1(\Delta X) = -(\Delta \Phi^\top \Delta X + \Delta X \Delta \Phi + \Psi^\top \Delta X \Delta \Psi + \Delta \Psi^\top \Delta X \Psi + \Delta \Psi^\top \Delta X \Delta \Psi), \quad (21b)$$

$$\begin{aligned}h_2(\Delta X) &= \tilde{\Psi}^\top \Delta X \tilde{D} \tilde{\Omega} \tilde{D}^\top \Delta X \tilde{\Psi} + \Delta X \tilde{B} \tilde{\Omega} \tilde{B}^\top \Delta X \\ &\quad + \Delta X \tilde{B} \tilde{\Omega} \tilde{D}^\top \Delta X \tilde{\Psi} + \tilde{\Psi}^\top \Delta X \tilde{D} \tilde{\Omega} \tilde{B}^\top \Delta X,\end{aligned} \quad (21c)$$

where  $\Omega = (I + \tilde{Q}(X)^{-1}\tilde{D}^\top \Delta X \tilde{D})^{-1}\tilde{Q}(X)^{-1}$ , the matrices  $\Delta A$ ,  $\Delta B$ , and so on are given by (9)-(12).

Note that  $E_1$  and  $E_2$  are not dependent on  $\Delta X$ ,  $h_1(\Delta X)$  is a linear function of  $\Delta X$ , and  $h_2(\Delta X)$  is a function of  $\Delta X$  with degree at most 2. Assume that the linear operator  $\mathcal{L}_c$  of (4) is invertible. It is easy to see that the perturbed equation (20) is true if and only if

$$\Delta X = \mathcal{L}_c^{-1}E_1 + \mathcal{L}_c^{-1}E_2 + \mathcal{L}_c^{-1}h_1(\Delta X) + \mathcal{L}_c^{-1}h_2(\Delta X). \quad (22)$$

Thus far, we have not specified the condition for the existence of the solution  $\Delta X$  in (22). In the subsequent discussion, we shall limit our attention to identifying the condition of the existence of a fixed point of (23), that is, to determine an upper bound on the size of  $\Delta X$ .

### 3 Perturbation bounds

Let  $f: \mathcal{S}^n \rightarrow \mathcal{S}^n$  be a continuous mapping defined by

$$f(Y) = \mathcal{L}_c^{-1}E_1 + \mathcal{L}_c^{-1}E_2 + \mathcal{L}_c^{-1}h_1(Y) + \mathcal{L}_c^{-1}h_2(Y) \quad \text{for } Y \in \mathcal{S}^n. \quad (23)$$

We see that any fixed point of the mapping  $f$  is a solution to the perturbed equation (22). Our approach in this section is to present an upper bound for the existence of some fixed points  $\Delta X$ . It starts with the discussion that the mapping  $f$  given by (23) satisfies

$$\|f(\Delta X)\|_F \leq \|\mathcal{L}_c^{-1}E_1\|_F + \|\mathcal{L}_c^{-1}E_2\|_F + \|\mathcal{L}_c^{-1}h_1(\Delta X)\|_F + \|\mathcal{L}_c^{-1}h_2(\Delta X)\|_F.$$

Define linear operators  $\mathcal{M}: \mathbb{R}^{n \times n} \rightarrow \mathcal{S}^n$ ,  $\mathcal{N}: \mathbb{R}^{n \times m} \rightarrow \mathcal{S}^n$ ,  $\mathcal{T}: \mathcal{S}^m \rightarrow \mathcal{S}^n$  and  $\mathcal{H}: \mathbb{R}^{n \times m} \rightarrow \mathcal{S}^n$  by

$$\mathcal{M}\Delta C = \mathcal{L}_c^{-1}(K^\top \Delta C + \Delta C^\top K), \quad (24a)$$

$$\mathcal{N}\Delta D = \mathcal{L}_c^{-1}(K^\top \Delta D F + F^\top \Delta D^\top K), \quad (24b)$$

$$\mathcal{T}\Delta R = \mathcal{L}_c^{-1}(F^\top \Delta R F), \quad (24c)$$

$$\mathcal{H}\Delta S = \mathcal{L}_c^{-1}(F^\top \Delta S^\top + \Delta S F), \quad (24d)$$

and the scalars  $\omega$ ,  $\mu$ ,  $\nu$ ,  $\tau$ ,  $\eta$  by

$$\omega = \|\mathcal{L}_c^{-1}\|, \quad \mu = \|\mathcal{M}\|, \quad \nu = \|\mathcal{N}\|, \quad \tau = \|\mathcal{T}\|, \quad \eta = \|\mathcal{H}\|. \quad (25)$$

From (21a) we then have

$$\|\mathcal{L}_c^{-1}E_1\|_F \leq \mu\|\Delta C\|_F + \nu\|\Delta D\|_F + \tau\|\Delta S\|_F + \eta\|\Delta R\|_F + \omega\|\Delta H\| \equiv \varepsilon_1. \quad (26)$$

We now move into more specific details pertaining to the discussion of the fixed point of the continuous mapping  $f$ . Before doing so, we need to describe an important property of the norm of the product of two matrices and repeatedly employ it in the following discussion. For the proof, the reader is referred to [30, Theorem 3.9].



**Lemma 3.1** *Let  $A$  and  $B$  be two matrices in  $\mathbb{R}^{n \times n}$ . Then  $\|AB\|_F \leq \|A\|_2 \|B\|_F$  and  $\|AB\|_F \leq \|A\|_F \|B\|_2$ .*

It immediately follows that the matrices  $\delta Q$  and  $\delta S$ , defined by (10), satisfy

$$\begin{aligned}\|\delta Q\|_F &\leq \|\Delta R\|_F + 2\|XD\|_2 \|\Delta D\|_F + \|X\|_2 \|\Delta D\|_F^2 \equiv \delta_r, \\ \|\delta S\|_F &\leq \|\Delta S\|_F + \|XC\|_2 \|\Delta D\|_F + \|XD\|_2 \|\Delta C\|_F \\ &\quad + \|X\|_2 \|\Delta D\|_F \|\Delta C\|_F + \|X\|_2 \|\Delta B\|_F \equiv \delta_s.\end{aligned}\quad (27)$$

Assume that the scalar  $\delta_r$  satisfies

$$1 - \|\mathcal{Q}(X)^{-1}\|_2 \delta_r > 0. \quad (28)$$

Then  $\|\mathcal{L}_c^{-1} E_2\|_F$  is bounded by

$$\begin{aligned}\|\mathcal{L}_c^{-1} E_2\|_F &\leq 2\omega \|X\|_2 (\|\Delta A\|_F + \|F\|_2 \|\Delta B\|_F) + \omega \|X\|_2 (\|\Delta C\|_F + \|F\|_2 \|\Delta D\|_F)^2 \\ &\quad + \frac{\omega \|\mathcal{Q}(X)^{-1}\|_2 (\|F\|_2 \delta_r + \delta_s)^2}{1 - \|\mathcal{Q}(X)^{-1}\|_2 \delta_r} \equiv \varepsilon_2.\end{aligned}\quad (29)$$

From (21b) we see that

$$\|h_1(\Delta X)\|_F \leq (2\|\Delta \Phi\|_F + 2\|\Psi\|_2 \|\Delta \Psi\|_F + \|\Delta \Psi\|_F^2) \|\Delta X\|_F,$$

and also from (18) and (19) we have

$$\begin{aligned}\|\Delta \Phi\|_F &\leq \|\Delta A\|_F + \|F\|_2 \|\Delta B\|_F + (\|B\mathcal{Q}(X)^{-1}\|_2 + \|\mathcal{Q}(X)^{-1}\|_2 \|\Delta B\|_F) \delta_s \\ &\quad + \frac{\|\mathcal{Q}(X)^{-1}\|_2 (\|B\|_2 + \|\Delta B\|_F) (\|F\|_2 + \|\mathcal{Q}(X)^{-1}\|_2 \delta_s) \delta_r}{1 - \|\mathcal{Q}(X)^{-1}\|_2 \delta_r} \equiv \delta_\Phi,\end{aligned}\quad (30)$$

$$\begin{aligned}\|\Delta \Psi\|_F &\leq \|\Delta C\|_F + \|F\|_2 \|\Delta D\|_F + (\|D\mathcal{Q}(X)^{-1}\|_2 + \|\mathcal{Q}(X)^{-1}\|_2 \|\Delta D\|_F) \delta_s \\ &\quad + \frac{\|\mathcal{Q}(X)^{-1}\|_2 (\|D\|_2 + \|\Delta D\|_F) (\|F\|_2 + \|\mathcal{Q}(X)^{-1}\|_2 \delta_s) \delta_r}{1 - \|\mathcal{Q}(X)^{-1}\|_2 \delta_r} \equiv \delta_\Psi.\end{aligned}\quad (31)$$

It follows that

$$\|\mathcal{L}_c^{-1} h_1(\Delta X)\|_F \leq \omega \delta \|\Delta X\|_F, \quad (32)$$

where the positive scalar  $\delta$  is defined by

$$\delta = 2\delta_\Phi + 2\psi \delta_\Psi + \delta_\Psi^2 \quad \text{with } \|\Psi\|_2 = \psi. \quad (33)$$

Also, from (28) and Lemma 3.1 we know that  $\tilde{\mathcal{Q}}(X) = \mathcal{Q}(X) + \delta Q = \mathcal{Q}(X)^{-1}(I + \mathcal{Q}(X)^{-1} \delta Q)$  and  $\|\mathcal{Q}(X)^{-1} \delta Q\|_F \leq \|\mathcal{Q}(X)^{-1}\|_2 \delta_r < 1$ . This implies that  $\tilde{\mathcal{Q}}(X)$  is nonsingular,

$$\begin{aligned}\|\tilde{\mathcal{Q}}\|_2^2 \|\tilde{\mathcal{Q}}(X)^{-1}\|_2 &= \|B + \Delta B\|_2^2 \|(\mathcal{Q}(X) + \delta Q)^{-1}\|_2 \\ &= \|B + \Delta B\|_2^2 \|(I + \mathcal{Q}(X)^{-1} \delta Q)^{-1} \mathcal{Q}(X)^{-1}\|_2 \\ &\leq \frac{\|\mathcal{Q}(X)^{-1}\|_2 (\|B\|_2 + \|\Delta B\|_F)^2}{1 - \|\mathcal{Q}(X)^{-1}\|_2 \delta_r} \equiv \gamma_B.\end{aligned}\quad (34)$$

Similarly, we have

$$\|\tilde{D}\|_2^2 \|\tilde{Q}(X)^{-1}\|_2 \leq \frac{\|\tilde{Q}(X)^{-1}\|_2 (\|D\|_2 + \|\Delta D\|_F)^2}{1 - \|\tilde{Q}(X)^{-1}\|_2 \delta_r} \equiv \gamma_D. \quad (35)$$

Assume that

$$1 - \gamma_D \|\Delta X\|_F > 0. \quad (36)$$

It then follows from Lemma 3.1 and (21c) that

$$\|h_2(\Delta X)\|_F \leq \frac{(\|\tilde{\Psi}\|_2 \|\tilde{D}\|_2 + \|\tilde{B}\|_2)^2 \|\tilde{Q}(X)^{-1}\|_2 \|\Delta X\|_F^2}{1 - \|\tilde{D}\|_2^2 \|\tilde{Q}(X)^{-1}\|_2 \|\Delta X\|_F}, \quad (37)$$

and from (9), (11) and (31) that

$$\begin{aligned} \|\tilde{B}\|_2 &\leq \|B\|_2 + \|\Delta B\|_2 \equiv \alpha_B, \\ \|\tilde{D}\|_2 &\leq \|D\|_2 + \|\Delta D\|_2 \equiv \alpha_D, \\ \|\tilde{\Psi}\|_2 &\leq \|\Psi\|_2 + \|\Delta \Psi\|_F \leq \|\Psi\|_2 + \delta_\Psi \equiv \tilde{\psi}. \end{aligned} \quad (38)$$

Upon substituting (34), (35) and (38) into (37), we see that

$$\|\mathcal{L}_c^{-1} h_2(\Delta X)\|_F \leq \frac{\omega(\tilde{\psi}^2 \gamma_D + 2\tilde{\psi} \alpha_B \alpha_D + \gamma_B) \|\Delta X\|_F^2}{1 - \gamma_D \|\Delta X\|_F}.$$

Finally, by (26), (29) and (32), we arrive at the statement

$$\|f(\Delta X)\|_F \leq \varepsilon + \omega \delta \|\Delta X\|_F + \frac{\omega \alpha \|\Delta X\|_F^2}{1 - \gamma_D \|\Delta X\|_F}, \quad (39)$$

where

$$\alpha \equiv \tilde{\psi}^2 \gamma_D + 2\tilde{\psi} \alpha_B \alpha_D + \gamma_B, \quad \varepsilon \equiv \varepsilon_1 + \varepsilon_2. \quad (40)$$

Consider the quadratic equation

$$(\gamma_D - \omega \delta \gamma_D + \omega \alpha) \xi^2 - (1 - \omega \delta + \varepsilon \gamma_D) \xi + \varepsilon = 0. \quad (41)$$

It is true that if

$$\delta < \frac{1}{\omega}, \quad (42a)$$

$$\varepsilon \leq \frac{(1 - \omega \delta)^2}{\gamma_D - \omega \delta \gamma_D + 2\omega \alpha + \sqrt{(\gamma_D - \omega \delta \gamma_D + 2\omega \alpha)^2 - \gamma_D^2 (1 - \omega \delta)^2}}, \quad (42b)$$

then the positive scalar  $\xi_*$  denoted by

$$\xi_* = \frac{2\varepsilon}{(1 - \omega \delta + \varepsilon \gamma_D) + \sqrt{(1 - \omega \delta + \varepsilon \gamma_D)^2 - 4(\gamma_D - \omega \delta \gamma_D + \omega \alpha) \varepsilon}} \quad (43)$$

is a solution to (41). Let  $\mathcal{S}_{\xi_*}^n$  be a compact subset of  $\mathcal{S}^n$  given by

$$\mathcal{S}_{\xi_*}^n = \{\Delta X \in \mathcal{S}^n : \|\Delta X\|_F \leq \xi_*\}.$$

It can be seen that in (39)

$$\|f(\Delta X)\|_F \leq \xi_* \quad \text{if } \Delta X \in \mathcal{S}_{\xi_*}^n.$$

It then follows from the Brouwer fixed-point theorem (see [31]) that the continuous mapping  $f$  has a fixed point  $\Delta X_* \in \mathcal{S}_{\xi_*}^n$ , that is, condition (22) automatically holds.

Observe also that if  $\Delta X \in \mathcal{S}_{\xi_*}^n$ , then

$$\begin{aligned} 1 - \gamma_D \|\Delta X\|_F &\geq (1 - \gamma_D \xi_*) \quad (\text{by (43)}) \\ &\geq 1 - \frac{2\varepsilon\gamma_D}{1 - \omega\delta + \varepsilon\gamma_D} = \frac{1 - \omega\delta - \varepsilon\gamma_D}{1 - \omega\delta + \varepsilon\gamma_D} \quad (\text{by (42b)}) \\ &\geq \frac{1 - \omega\delta - \frac{(1-\omega\delta)^2\gamma_D}{\gamma_D - \omega\delta\gamma_D + 2\omega\alpha}}{1 - \omega\delta + \varepsilon\gamma_D} \quad (\text{by (42a)}) \\ &= \frac{2(1 - \omega\delta)\omega\alpha}{(1 - \omega\delta + \varepsilon\gamma_D)(\gamma_D - \omega\delta\gamma_D + 2\omega\alpha)} \geq 0. \end{aligned}$$

This implies that assumption (36) is true, if assumption (42a)-(42b) is true.

#### 4 Stability analysis

We have shown that the mapping  $f$  given by (23) has a Hermitian fixed point  $\Delta X_*$ . This further implies that perturbed SARE (1a)-(1b) has a Hermitian solution  $\tilde{X} = X + \Delta X_*$ . In this section, we want to discuss the stability of the solution  $\tilde{X}$ , i.e., show that the solution  $\tilde{X}$  is the unique maximal solution to SARE (1a)-(1b). Let  $\Upsilon$  and  $\Pi$  be two operators defined by

$$\Upsilon(W) = \Phi^\top W + W\Phi, \quad \Pi(W) = \Psi^\top W\Psi, \quad W \in \mathcal{S}^n,$$

with the notations  $\Phi$  and  $\Psi$  given in Definition 1.1. It follows that the operator  $\mathcal{L}_c$  defined by (4) can also be written as

$$\mathcal{L}_c(W) = \Upsilon(W) + \Pi(W), \quad W \in \mathcal{S}^n. \quad (44)$$

We then have the following important result addressing the condition for a linear operator to be stable. To see a few necessary and sufficient conditions on the stability, we refer to the results and proofs given in [3].

**Theorem 4.1** *The linear operator  $\mathcal{L}_c = \Upsilon + \Pi$  given by (44) is stable, i.e.,  $\sigma(\mathcal{L}_c) \subset \mathbb{C}_-$ , if and only if  $\sigma(\Phi) \subset \mathbb{C}_-$  and  $\det(\Upsilon + \tau\Pi) \neq 0$  for all  $\tau \in [0, 1]$ .*

When small perturbations  $Z_1, Z_2 \in \mathbb{R}^{n \times n}$  are taken into consideration, the perturbed operator of  $\mathcal{L}_c$  can be expressed by

$$\tilde{\mathcal{L}}_c(W) = \tilde{\Upsilon}(W) + \tilde{\Pi}(W), \quad (45)$$

where  $\tilde{\Upsilon}(W) = (\Phi + Z_1)^\top W + W(\Phi + Z_1)$  and  $\tilde{\Pi}(W) = (\Psi + Z_2)^\top W + W(\Psi + Z_2)$  for all  $W \in \mathcal{S}^n$ . Define the quantity

$$\ell(\theta) = \|(\Upsilon + \theta\Pi)^{-1}\|$$

for  $\theta \in [0, 1]$  and

$$\beta(\mathcal{L}_c) = \min_{(Z_1, Z_2) \in \mathcal{Z}} \max\{\|Z_1\|, \|Z_2\|\},$$

where the set  $\mathcal{Z} = \{(Z_1, Z_2) \in \mathbb{R}^{n \times n} \times \mathbb{R}^{n \times n} \mid \det(\tilde{\Upsilon} + \theta\tilde{\Pi}) = 0 \text{ for some } \theta \in [0, 1]\}$ . It should be noted that if  $\sigma(\Phi) \subset \mathbb{C}_-$ ,  $\Psi = 0$  and  $Z_2 = 0$ , then

$$\beta(\mathcal{L}_c) \leq \beta(\Phi), \quad (46)$$

where the value  $\beta(\Phi)$  is defined by [27]

$$\beta(\Phi) = \min\left\{\|Z_1\| \mid \max_{1 \leq j \leq n} \operatorname{Re} \lambda_j(\Phi + Z_1) = 0, Z_1 \in \mathbb{R}^{n \times n}\right\}. \quad (47)$$

Here,  $\lambda_j(\Phi + Z_1)$  ( $j = 1, \dots, n$ ) denote the eigenvalues of  $\Phi + Z_1$ .

The connection between  $\beta(\mathcal{L}_c)$  and the maximum of the scalar function  $\ell(\theta)$  on  $[0, 1]$  can be established in the following form.

**Theorem 4.2** [23] *Suppose that the linear operator  $\mathcal{L}_c$  given by (44) is stable, and let*

$$\ell_c = \max_{\theta \in [0, 1]} \ell(\theta), \quad \psi = \|\Psi\|. \quad (48)$$

*Then*

$$\beta(\mathcal{L}_c) \geq \frac{\ell_c^{-1}}{(\psi + 1) + \sqrt{(\psi + 1)^2 + \ell_c^{-1}}}.$$

We now apply Theorem 4.2 to (46) and obtain that

$$\beta(\Phi) \geq \beta(\mathcal{L}_c) \geq \frac{\ell_c^{-1}}{(\psi + 1) + \sqrt{(\psi + 1)^2 + \ell_c^{-1}}}.$$

Hence, if a perturbation matrix  $Z_1 \in \mathbb{R}^{n \times n}$  satisfies

$$\|Z_1\| < \frac{\ell_c^{-1}}{(\psi + 1) + \sqrt{(\psi + 1)^2 + \ell_c^{-1}}},$$

then (47) implies that the matrix  $\Phi + Z_1$  must be  $c$ -stable.

We now turn to a key stability test of the operator  $\mathcal{L}_c$ , the striking tool of our stability analysis.

**Theorem 4.3** [23] *Suppose that the linear operator  $\mathcal{L}_c$  is stable, and let the scalars  $\ell_c$  and  $\psi$  be defined as in (48). If the perturbation matrices  $Z_1, Z_2 \in \mathbb{R}^{n \times n}$  satisfy*

$$\max\{\|Z_1\|, \|Z_2\|\} < \frac{\ell_c^{-1}}{(\psi + 1) + \sqrt{(\psi + 1)^2 + \ell_c^{-1}}}, \quad (49)$$

then  $\Phi + Z_1$  is  $c$ -stable and the perturbed linear operator  $\tilde{\mathcal{L}}_c$  defined by (45) is also stable, i.e.,  $\sigma(\tilde{\mathcal{L}}_c) \subset \mathbb{C}_-$ .

Upon substituting  $\tilde{X}$  for  $X$  in  $\tilde{S}(X)$  and  $\tilde{Q}(X)$  of (11), we shall have

$$\tilde{Q}(\tilde{X}) = \tilde{R} + \tilde{D}^\top \tilde{X} \tilde{D} = Q(X) + \delta Q + \tilde{D}^\top \Delta X \tilde{D} \equiv Q(X) + \Delta R, \quad (50)$$

$$\begin{aligned} \tilde{S}(\tilde{X}) &= \tilde{S} + \tilde{C}^\top \tilde{X} \tilde{D} + \tilde{X} \tilde{B} = S(X) + \delta S + \Delta X \tilde{B} + \tilde{C}^\top \Delta X \tilde{D} \\ &\equiv S(X) + \Delta S. \end{aligned} \quad (51)$$

Also, corresponding to  $\tilde{X}$ , the perturbed  $\Psi_{\tilde{X}}$  and  $\Phi_{\tilde{X}}$  of  $\Psi$  and  $\Phi$ , respectively, can be expressed in terms of the formulae

$$\begin{aligned} \Phi_{\tilde{X}} &= (\Delta A + A) - (\Delta B + B)(\Gamma + \Delta S)^{-1}(\Lambda + \Delta R)^\top := \Phi + \Delta \Phi, \\ \Psi_{\tilde{X}} &= (\Delta C + C) - (\Delta D + D)(\Gamma + \Delta R)^{-1}(\Lambda + \Delta S)^\top := \Psi + \Delta \Psi, \end{aligned} \quad (52)$$

with

$$\begin{aligned} \Delta \Phi &:= \Delta A - \Delta B Q(X)^{-1} S(X)^\top - \Delta B Q(X)^{-1} \Delta S^\top - B Q(X)^{-1} \Delta S^\top \\ &\quad + (\Delta B + B) Q(X)^{-1} \Delta R Q(X)^{-1} (I + Q(X)^{-1} \Delta R)^{-1} (S(X)^\top + \Delta S^\top), \\ \Delta \Psi &:= \Delta C - \Delta D Q(X)^{-1} S(X)^\top - \Delta D Q(X)^{-1} \Delta S^\top - D Q(X)^{-1} \Delta S^\top \\ &\quad + (\Delta D + D) Q(X)^{-1} \Delta R Q(X)^{-1} (I + Q(X)^{-1} \Delta R)^{-1} (S(X)^\top + \Delta S^\top). \end{aligned}$$

Let  $\alpha_C := \|C\|_2 + \|\Delta C\|_2$ . Since  $\|\Delta X\|_F \leq \xi_*$ , it follows from (38), (50) and (51) that

$$\begin{aligned} \|\Delta R\|_F &\leq \delta_r + \alpha_D^2 \xi_* := c_r, \\ \|\Delta S\|_F &\leq \delta_s + (\alpha_B + \alpha_C \alpha_D) \xi_* := c_s. \end{aligned}$$

Thus  $\|\Delta \Phi\|_F$  and  $\|\Delta \Psi\|_F$  are bounded by the inequalities

$$\begin{aligned} \|\Delta \Phi\|_F &\leq \|\Delta A\|_F + \|F\|_2 \|\Delta B\|_F + \frac{\alpha_B \|Q(X)^{-1}\|_2 (c_s + c_r \|F\|_2)}{1 - c_r \|Q(X)^{-1}\|_2}, \\ \|\Delta \Psi\|_F &\leq \|\Delta C\|_F + \|F\|_2 \|\Delta D\|_F + \frac{\alpha_D \|Q(X)^{-1}\|_2 (c_s + c_r \|F\|_2)}{1 - c_r \|Q(X)^{-1}\|_2}. \end{aligned}$$

Here, the above upper bounds are obtained by simplifying those given by (30) and (31). Let

$$\begin{aligned} f &= \|F\|_2, \quad \gamma = \|Q(X)^{-1}\|_2, \\ \alpha_C &= \|C\|_2 + \|\Delta C\|_F, \quad \zeta_1 = \max\{\|\Delta A\|_F, \|\Delta C\|_F\}, \\ \zeta_2 &= \max\{\|\Delta B\|_F, \|\Delta D\|_F\}, \quad \zeta_3 = \max\{\alpha_B, \alpha_D\}, \end{aligned} \quad (53)$$

where  $\alpha_B$  and  $\alpha_D$  are defined by (38) and  $\Theta$  is defined to be the right-hand side of (49), that is,

$$\Theta = \frac{\ell_c^{-1}}{(\psi + 1) + \sqrt{(\psi + 1)^2 + \ell_c^{-1}}}. \quad (54)$$

We then have

$$\max\{\|\Delta\Phi\|_F, \|\Delta\Psi\|_F\} \leq \zeta_1 + f\zeta_2 + \frac{\gamma\zeta_3(\delta_s + \delta_r f) + \gamma\zeta_3(\alpha_B + \alpha_C\alpha_D + \alpha_D^2 f)\xi_*}{1 - \delta_r f - \alpha_D^2 f\xi_*}.$$

It follows that if the condition

$$\zeta_1 + f\zeta_2 + \frac{\gamma\zeta_3(\delta_s + \delta_r f) + \gamma\zeta_3(\alpha_B + \alpha_C\alpha_D + \alpha_D^2 f)\xi_*}{1 - \delta_r f - \alpha_D^2 f\xi_*} < \Theta$$

or, equivalently,

$$\xi_* < \frac{(\Theta - \zeta_1 - f\zeta_2)(1 - \delta_r f) - \gamma\zeta_3(\delta_s + \delta_r f)}{(\Theta - \zeta_1 - f\zeta_2)\alpha_D^2 f + \gamma\zeta_3(\alpha_B + \alpha_C\alpha_D + \alpha_D^2 f)}$$

holds, then corresponding to Theorem 4.3, the perturbed linear operator  $\tilde{\mathcal{L}}_c$  with respect to  $\tilde{X}$  is stable. In other words, the matrix  $\tilde{X} \in \mathcal{S}^n$  must be the unique stabilizing (and maximal) solution to perturbed SARE (5a)-(5b).

We now have all the materials needed for the existence of a stabilizing solution of (5a)-(5b).

**Theorem 4.4** (Perturbation bound) *Let  $X$  be the stabilizing solution of (1a)-(1b). Let  $\omega, \delta_r, \delta_s, \delta, \gamma_D, \alpha_B, \alpha_D, \alpha, \varepsilon, f, \gamma, \alpha_C, \zeta_1, \zeta_2, \zeta_3, \Theta$  be the scalars defined by (25), (27), (33), (35), (38), (40), (53) and (54), respectively. Define*

$$\xi_* = \frac{2\varepsilon}{(1 - \omega\delta + \varepsilon\gamma_D) + \sqrt{(1 - \omega\delta + \varepsilon\gamma_D)^2 - 4(\gamma_D - \omega\delta\gamma_D + \omega\alpha)\varepsilon}}.$$

*If the perturbed quantities of the coefficients of (5a)-(5b) are sufficiently small, for example,  $\varepsilon \ll 1$ , such that*

$$1 - \|\mathcal{Q}(X)^{-1}\|_2 \delta_r > 0,$$

$$1 - \omega\delta > 0,$$

$$\frac{(1 - \omega\delta)^2}{\gamma_D - \omega\delta\gamma_D + 2\omega\alpha + \sqrt{(\gamma_D - \omega\delta\gamma_D + 2\omega\alpha)^2 - \gamma_D^2(1 - \omega\delta)^2}} - \varepsilon \geq 0,$$

$$\frac{(\Theta - \zeta_1 - f\zeta_2)(1 - \delta_r f) - \gamma\zeta_3(\delta_s + \delta_r f)}{(\Theta - \zeta_1 - f\zeta_2)\alpha_D^2 f + \gamma\zeta_3(\alpha_B + \alpha_C\alpha_D + \alpha_D^2 f)} - \xi_* > 0,$$

*then perturbed SARE (5a)-(5b) has the unique stabilizing solution  $\tilde{X}$ , and*

$$\frac{\|\tilde{X} - X\|_F}{\|X\|_F} \leq \frac{\xi_*}{\|X\|_F}. \quad (55)$$

## 5 Condition number of the SARE

In the study of a computational problem, a fundamental issue is to determine the condition number of a problem to be the ratio of the relative change in the solution to the relative

change in the argument. Applying the theory of condition number given by Rice [32], we define the condition number  $c(X)$  of the stabilizing solution  $X$  of SARE (1a)-(1b) by

$$c(X) = \lim_{\delta \rightarrow 0^+} \sup_{\Omega_\delta} \frac{\|\Delta X\|_F}{\kappa \delta}, \quad (56)$$

where the set of perturbed matrices  $\Omega_\delta$  is defined by

$$\begin{aligned} \Omega_\delta &= \Omega_\delta(\kappa_A, \kappa_B, \kappa_C, \kappa_D, \kappa_S, \kappa_R, \kappa_H) \\ &= \left\{ (\Delta A, \Delta B, \Delta C, \Delta D, \Delta S) \in \mathbb{R}^{n \times n} \times \mathbb{R}^{n \times m} \times \mathbb{R}^{n \times n} \times \mathbb{R}^{n \times m} \times \mathbb{R}^{n \times m}, \right. \\ &\quad \left. (\Delta R, \Delta H) \in \mathcal{S}^n \times \mathcal{S}^m \mid 0 < \delta_p \leq \delta \right\}, \end{aligned} \quad (57)$$

with

$$\delta_p = \left\| \left( \frac{\Delta A}{\kappa_A}, \frac{\Delta B}{\kappa_B}, \frac{\Delta C}{\kappa_C}, \frac{\Delta D}{\kappa_D}, \frac{\Delta S}{\kappa_S}, \frac{\Delta R}{\kappa_R}, \frac{\Delta H}{\kappa_H} \right) \right\|_F,$$

and  $\kappa_A, \kappa_B, \kappa_C, \kappa_D, \kappa_S, \kappa_R, \kappa_H, \kappa$  are positive parameters. Then (56) gives the absolute condition number  $c_{\text{abs}}(X)$  if

$$(\kappa_A, \kappa_B, \kappa_C, \kappa_D, \kappa_S, \kappa_R, \kappa_H, \kappa) = (1, 1, 1, 1, 1, 1, 1, 1)$$

and gives the relative condition number  $c_{\text{rel}}(X)$  if

$$(\kappa_A, \kappa_B, \kappa_C, \kappa_D, \kappa_S, \kappa_R, \kappa_H, \kappa) = (\|A\|_F, \|B\|_F, \|C\|_F, \|D\|_F, \|S\|_F, \|R\|_F, \|H\|_F, \|X\|_F).$$

It follows from (22) and (24a)-(24d) that

$$\Delta X = \mathcal{P} \Delta A + \mathcal{Q} \Delta B + \mathcal{M} \Delta C + \mathcal{N} \Delta D + \mathcal{H} \Delta S + \mathcal{T} \Delta R + \mathcal{L}_c^{-1} \Delta H + O(\delta_p^2),$$

where the linear operators  $\mathcal{P} : \mathbb{R}^{n \times n} \rightarrow \mathcal{S}^n$  and  $\mathcal{Q} : \mathbb{R}^{n \times m} \rightarrow \mathcal{S}^n$  are defined by

$$\mathcal{P} \Delta A = \mathcal{L}_c^{-1} (X \Delta A + \Delta A^\top X), \quad (58a)$$

$$\mathcal{Q} \Delta B = \mathcal{L}_c^{-1} (X \Delta B F + F^\top \Delta B^\top X). \quad (58b)$$

In order to derive the explicit expression for the condition number  $c(X)$  of the stabilizing solution  $X$  of (1a)-(1b), we require a theorem concerning the form of the optimal solution. This theorem can be regarded as a theoretical extension of the results discussed in [25, 33]. Most strategies have been established earlier by using much heavier machinery. Since this theorem is most relevant to our stability analysis, we briefly outline a direct proof with ideas from [34] to make this presentation more self-contained.

**Theorem 5.1** *Let  $\mathcal{L} : \mathbb{R}^{n \times n} \times \mathbb{R}^{m \times m} \rightarrow \mathbb{R}^{k \times k}$  be a linear operator and*

$$\mathcal{L}(Z_1, Z_2)^\top = \mathcal{L}(Z_1, Z_2^\top) \quad (59)$$

for all  $Z_1 \in \mathbb{R}^{n \times n}$  and  $Z_2 \in \mathbb{R}^{m \times m}$ . Then the optimal solution  $(Z_{1*}, Z_{2*})$  to the problem

$$\max_{\|(Z_1, Z_2)\|_F=1} \|\mathcal{L}(Z_1, Z_2)\|_F \quad (60)$$

exists for some  $Z_{1*} \in \mathbb{R}^{n \times n}$  and  $Z_{2*} = \pm Z_{2*}^\top \in \mathbb{R}^{m \times m}$ . Furthermore, if the linear operator  $\mathcal{L}(0, Z_2)$  is a positive operator with respect to any  $Z_2 \in \mathbb{C}^{m \times m}$ , that is, for all  $Z_2 \in \mathbb{C}^{m \times m}$ , we have

$$\mathcal{L}(0, Z_2) \succeq 0 \quad \text{if } Z_2 \succeq 0.$$

Then there exists an optimal solution  $(Z_{1*}, Z_{2*}) \in \mathbb{R}^{n \times n} \times \mathbb{R}^{m \times m}$  to problem (60) such that  $Z_{2*}$  is symmetric.

*Proof* Since  $\mathcal{L}$  is a linear operator on a finite dimensional space, it is clear that the optimal solution of (60) exists. Assume that  $(Z_{1*}, Z_{2*})$  solves this optimization problem. Let  $\sigma_{\max} = \|L(Z_{1*}, Z_{2*})\|_F$  and  $L \in \mathbb{R}^{k^2 \times (n^2 + m^2)}$  be the matrix representation of the operator  $\mathcal{L}$  such that

$$\text{vec}(\mathcal{L}(Z_1, Z_2)) = L \begin{bmatrix} \text{vec}(Z_1) \\ \text{vec}(Z_2) \end{bmatrix}. \quad (61)$$

By (59) and (61), we have

$$\begin{bmatrix} \text{vec}(Z_{1*}) \\ \text{vec}(Z_{2*}) \end{bmatrix}^\top L^\top L \begin{bmatrix} \text{vec}(Z_{1*}) \\ \text{vec}(Z_{2*}) \end{bmatrix} = \begin{bmatrix} \text{vec}(Z_{1*}) \\ \text{vec}(Z_{2*}^\top) \end{bmatrix}^\top L^\top L \begin{bmatrix} \text{vec}(Z_{1*}) \\ \text{vec}(Z_{2*}^\top) \end{bmatrix} = \sigma_{\max}^2. \quad (62)$$

Note that

$$\ell := \|Z_{1*}\|_F^2 + \left\| \frac{1}{2}(Z_{2*} + Z_{2*}^\top) \right\|_F^2 = 0, \quad \text{only if } Z_{1*} = 0, Z_{2*} = -Z_{2*}^\top.$$

It follows that if  $Z_{2*}^\top \neq -Z_{2*}$ , by (62), we see that  $\frac{1}{\sqrt{\ell}}(Z_{1*}, \frac{1}{2}(Z_{2*} + Z_{2*}^\top))$  is another optimal solution for (60). This proves the first part of the theorem.

For the second part, if there exists a symmetric optimal solution, then it completes the proof. Otherwise, from the first part, we know that there exists an optimal solution  $(Z_{1*}, Z_{2*})$  with  $Z_{1*} = 0$ ,  $Z_{2*} = -Z_{2*}^\top \in \mathbb{R}^{m \times m}$  and  $\|(Z_{1*}, Z_{2*})\|_F = 1$  to (60). Let  $i = \sqrt{-1}$ . We have the following matrix decomposition:

$$Z_{2*} = Q^\top \text{diag} \left( \begin{bmatrix} 0 & -\omega_1 \\ \omega_1 & 0 \end{bmatrix}, \dots, \begin{bmatrix} 0 & -\omega_k \\ \omega_k & 0 \end{bmatrix}, 0_r \right) Q,$$

where  $0_r$  is a zero matrix with size  $r \times r$ ,  $Q$  is an  $m \times m$  orthogonal matrix, and  $\omega_j > 0$  for  $1 \leq j \leq k$ . Let

$$\tilde{Z}_{2*} := Q^\top \text{diag} \left( \begin{bmatrix} \omega_1 & 0 \\ 0 & \omega_1 \end{bmatrix}, \dots, \begin{bmatrix} \omega_k & 0 \\ 0 & \omega_k \end{bmatrix}, 0_r \right) Q$$



be a real symmetric matrix. Since  $-\tilde{Z}_{2*} \preceq iZ_{2*} \preceq \tilde{Z}_{2*}$ , it is true that

$$\mathcal{L}(0, \tilde{Z}_{2*}) \succeq \mathcal{L}(0, iZ_{2*}) \succeq \mathcal{L}(0, -\tilde{Z}_{2*}) = -\mathcal{L}(0, \tilde{Z}_{2*}).$$

Using the fact that  $\|iZ_{2*}\|_F = \|\tilde{Z}_{2*}\|_F$ , we see that  $\|(0, \tilde{Z}_{2*})\|_F = \|(0, iZ_{2*})\|_F = \|(0, Z_{2*})\|_F = 1$  and

$$\|\mathcal{L}(0, \tilde{Z}_{2*})\|_F \geq \|\mathcal{L}(0, iZ_{2*})\|_F = \|\mathcal{L}(0, Z_{2*})\|_F.$$

If  $W_1 \succeq W_2 \succeq -W_1$ , then  $\|W_1\|_F \geq \|W_2\|_F$ , which implies that  $(0, \tilde{Z}_{2*})$  is a symmetric optimal solution to (60) (see [25, Lemma A.1]). This completes the proof.  $\square$

With the existence theory established above, it is interesting to note that the condition number  $c(X)$  defined by (56) can be written as

$$\begin{aligned} c(X) &= \frac{1}{\kappa} \lim_{\delta \rightarrow 0^+} \sup_{\Omega_\delta} \frac{\|\mathcal{P}\Delta A + \mathcal{Q}\Delta B + \mathcal{M}\Delta C + \mathcal{N}\Delta D + \mathcal{H}\Delta S + \mathcal{T}\Delta R + \mathcal{L}_c^{-1}\Delta H\|_F}{\delta} \\ &= \frac{1}{\kappa} \max_{\delta_p > 0} \frac{\|\mathcal{P}\Delta A + \mathcal{Q}\Delta B + \mathcal{M}\Delta C + \mathcal{N}\Delta D + \mathcal{H}\Delta S + \mathcal{T}\Delta R + \mathcal{L}_c^{-1}\Delta H\|_F}{\delta_p}. \end{aligned} \quad (63)$$

Note that the second equality in (63) is only an application of linearity of the norm. (For the proof, see Lemma A.1.) Observe further that the inverse operator  $\mathcal{L}_c^{-1}$  of (4) satisfies

$$[\mathcal{L}_c^{-1}(W)]^\top = \mathcal{L}_c^{-1}(W^\top)$$

since  $[\mathcal{L}_c(W)]^\top = \mathcal{L}_c(W^\top)$  for all  $W \in \mathbb{C}^{n \times n}$ . It follows that

$$\begin{aligned} [\mathcal{T}\Delta R]^\top &= \mathcal{T}\Delta R^\top, & [\mathcal{P}\Delta A]^\top &= \mathcal{P}\Delta A, & [\mathcal{Q}\Delta B]^\top &= \mathcal{Q}\Delta B, \\ [\mathcal{M}\Delta C]^\top &= \mathcal{M}\Delta C, & [\mathcal{N}\Delta D]^\top &= \mathcal{N}\Delta D, & [\mathcal{H}\Delta S]^\top &= \mathcal{H}\Delta S. \end{aligned}$$

Also, it is known that the inverse operator  $\mathcal{L}_c^{-1}$  is positive [3, Corollary 3.8]. It follows that  $\mathcal{T}$  is also a positive operator. Now, applying Theorem 5.1 to the operator  $\mathcal{P}\Delta A + \mathcal{Q}\Delta B + \mathcal{M}\Delta C + \mathcal{N}\Delta D + \mathcal{H}\Delta S + \mathcal{T}\Delta R + \mathcal{L}_c^{-1}\Delta H$  in (63), we obtain the equality

$$\begin{aligned} c(X) &= \frac{1}{\kappa} \max_{\Omega} \frac{\|\kappa_A \mathcal{P}\Delta A + \kappa_B \mathcal{Q}\Delta B + \kappa_C \mathcal{M}\Delta C + \kappa_D \mathcal{N}\Delta D + \kappa_S \mathcal{H}\Delta S + \kappa_R \mathcal{T}\Delta R + \kappa_H \mathcal{L}_c^{-1}\Delta H\|_F}{\|(\Delta A, \Delta B, \Delta C, \Delta D, \Delta S, \Delta R, \Delta H)\|_F}, \end{aligned}$$

where the extended set  $\tilde{\Omega}$  is defined by

$$\begin{aligned} \tilde{\Omega} &= \{(\Delta A, \Delta B, \Delta C, \Delta D, \Delta S, \Delta R, \Delta H) \in \mathbb{R}^{n \times n} \times \mathbb{R}^{n \times m} \times \mathbb{R}^{n \times n} \times \mathbb{R}^{n \times m} \\ &\quad \times \mathbb{R}^{n \times n} \times \mathbb{R}^{n \times n} \times \mathbb{R}^{m \times m} \mid \|(\Delta A, \Delta B, \Delta C, \Delta D, \Delta S, \Delta R, \Delta H)\|_F > 0\}. \end{aligned}$$

On the other hand, observe that the matrix representation of the operation  $\mathcal{L}_c$  in (4) can be written in terms of  $L_c = I \otimes \Phi + \Phi^\top \otimes I + \Psi^\top \otimes \Psi$ . Corresponding to (24a)-(24d) and

(58a)-(58b), we let

$$\begin{aligned}M_c &= L_c^{-1}(I \otimes K^\top + (K^\top \otimes I)P_{n,n}), \\N_c &= L_c^{-1}(F^\top \otimes K^\top + (K^\top \otimes F^\top)P_{m,n}), \\T_c &= L_c^{-1}(F^\top \otimes F^\top), \\H_c &= L_c^{-1}(F^\top \otimes I + (I \otimes F^\top)P_{m,n}), \\P_c &= L_c^{-1}(I \otimes X + (X \otimes I)P_{n,n}), \\Q_c &= L_c^{-1}(F^\top \otimes X + (X \otimes F^\top)P_{m,n})\end{aligned}$$

and

$$\mathcal{U} = (\kappa_A P_c, \kappa_B Q_c, \kappa_C M_c, \kappa_D N_c, \kappa_S H_c, \kappa_R T_c, \kappa_H L_c^{-1}).$$

It follows that

$$c(X) = \frac{1}{\kappa} \max_{V \in \tilde{\Omega}} \frac{\|\mathcal{U} \text{vec}(V)\|_2}{\|\text{vec}(V)\|_2} = \frac{\|\mathcal{U}\|_2}{\kappa}.$$

Based on the above discussion, we have the following result.

**Theorem 5.2** *The condition number  $c(X)$  given by (56) has the explicit expression  $\frac{\|\mathcal{U}\|_2}{\kappa}$ . In particular, we have the relative condition number*

$$c_{\text{rel}}(X) = \frac{\|(\|A\|_F P_c, \|B\|_F Q_c, \|C\|_F M_c, \|D\|_F N_c, \|S\|_F H_c, \|R\|_F T_c, \|H\|_F L_c^{-1})\|_2}{\|X\|_F}. \quad (64)$$

## 6 Numerical experiment

In this section we want to demonstrate the sharpness of perturbation bound (55) and its relationship with the relative condition number (64). Based on Newton's iteration [3], a numerical example, done with  $2 \times 2$  coefficient matrices, is illustrated. The numerical algorithm is described in Algorithm 1. The corresponding stopping criterion is determined when the value of the *Normalized Residual* (NRes)

$$\text{NRes} = \frac{\|\tilde{\mathcal{P}}(\tilde{X}) - \tilde{\mathcal{S}}(\tilde{X})\tilde{\mathcal{Q}}(\tilde{X})^{-1}\tilde{\mathcal{S}}(\tilde{X})^\top\|}{\|\tilde{\mathcal{P}}(\tilde{X})\| + \|\tilde{\mathcal{S}}(\tilde{X})\| \|\tilde{\mathcal{Q}}(\tilde{X})^{-1}\| \|\tilde{\mathcal{S}}(\tilde{X})^\top\|}$$

is less than or equal to a prescribed tolerance.

**Example 1** Given a parameter  $r = 10^{-m}$ , for some  $m > 0$ , let the matrices  $A, B, C, D$  be defined by

$$A = \begin{bmatrix} -4 & 0 \\ 0 & -1 \end{bmatrix}, \quad B = \begin{bmatrix} 1 & 0 \\ 0 & -\sqrt{r} \end{bmatrix}, \quad C = I_2, \quad D = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix},$$

and the matrices  $S, R, H$  be defined by

$$S = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}, \quad R = \begin{bmatrix} 3 & 0 \\ 0 & r \end{bmatrix}, \quad H = \begin{bmatrix} -13/2 & 0 \\ 0 & 1 \end{bmatrix}.$$

---

<b>Algorithm 1:</b> SARE	$[\tilde{X}] = \text{SARE}(\tilde{A}, \tilde{B}, \tilde{C}, \tilde{D}, \tilde{R}, \tilde{S}, \tilde{H})$
<hr/>	
<b>Input:</b> Matrices $\tilde{A}, \tilde{C} \in \mathbb{R}^{n \times n}$ , $\tilde{B}, \tilde{D}, \tilde{S} \in \mathbb{R}^{n \times m}$ , $\tilde{H} \in \mathbb{S}^n$ , $\tilde{R} \in \mathbb{S}^m$	
<b>Output:</b> Matrix $\tilde{X} \in \mathbb{S}^n$	
<b>begin</b>	
Choose $\tilde{X}_0 \in \mathbb{R}^{n \times n}$ ;	
<b>for</b> $i \leftarrow 1, \dots$ <b>do</b>	
$\tilde{F}_k = -(\tilde{R} + \tilde{D}^\top \tilde{X}_k \tilde{D})^{-1}(\tilde{B}^\top \tilde{X}_k \tilde{C} + \tilde{S}^\top)$ ; $\tilde{\Phi}_k = \tilde{A} + \tilde{B} \tilde{F}_k$ ; $\tilde{\Psi}_k = \tilde{C} + \tilde{D} \tilde{F}_k$ ;	
$\tilde{\mathcal{R}}(\tilde{X}_k)$ as defined by (7);	
Updating $\tilde{X}_{k+1}$ by solving	
$\tilde{\Phi}_k^\top \tilde{X}_{k+1} + \tilde{X}_{k+1} \tilde{\Phi}_k + \tilde{\Psi}_k^\top \tilde{X}_{k+1} \tilde{\Psi}_k = \tilde{\Phi}_k^\top \tilde{X}_k + \tilde{X}_k \tilde{\Phi}_k + \tilde{\Psi}_k^\top \tilde{X}_k \tilde{\Psi}_k - \tilde{\mathcal{R}}(\tilde{X}_k)$ ;	
<b>end</b>	
<b>end</b>	

---

**Table 1** Relative errors and perturbation bounds

$j$	Relative error	$\frac{\tilde{\epsilon}_*}{\ \tilde{X}\ _F}$
5	$6.28 \times 10^{-5}$	$5.74 \times 10^{-4}$
6	$6.34 \times 10^{-6}$	$5.22 \times 10^{-5}$
7	$1.14 \times 10^{-6}$	$8.42 \times 10^{-6}$
8	$1.61 \times 10^{-7}$	$7.42 \times 10^{-7}$
9	$1.05 \times 10^{-8}$	$5.58 \times 10^{-8}$

It is easily seen that the unique stabilizing and maximal solution is

$$X = \begin{bmatrix} -1 & 0 \\ 0 & (\sqrt{5} - 1)/2 \end{bmatrix}.$$

Let the perturbed coefficient matrices  $\Delta A$ ,  $\Delta B$ ,  $\Delta C$ ,  $\Delta D$ ,  $\Delta S$ ,  $\Delta R$  and  $\Delta H$  be generated using the MATLAB command randn with the weighted coefficient  $10^{-j}$ . That is, the matrices  $\Delta A$ ,  $\Delta B$ ,  $\Delta C$ ,  $\Delta D$ ,  $\Delta S$ ,  $\Delta R$  and  $\Delta H$  are generated in forms of  $\text{randn}(2) \times 10^{-j}$ , respectively. Since  $\Delta R$  and  $\Delta H$  are required to be symmetric, we need to fine-tune the perturbed matrices  $\Delta R$  and  $\Delta H$  by redefining  $\Delta R$  and  $\Delta H$  as  $\Delta R + \Delta R^\top$  and  $\Delta H + \Delta H^\top$ , respectively. Now, let  $(\tilde{A}, \tilde{B}, \tilde{C}, \tilde{D}, \tilde{S}, \tilde{R}, \tilde{H}) = (A + \Delta A, B + \Delta B, C + \Delta C, D + \Delta D, S + \Delta S, R + \Delta R, H + \Delta H)$ , which are coefficient matrices of SARE (5a)-(5b).

Firstly, we would like to evaluate the accuracy of the perturbation bound with the fixed parameter  $r = 10^{-2}$ , i.e.,  $m = 2$ , and different weighted coefficients,  $10^{-j}$ , for  $j = 5, \dots, 9$ . It can be seen from Table 1 that the values of the relative errors are closely bounded by our perturbation bounds of (55). In other words, (55) does provide a sharp upper bound of the relative errors of the stabilizing solution  $X$ .

Secondly, we want to investigate how ill-conditioned matrices affect the quantities of perturbation bounds. In this sense, the weighted coefficients are fixed to be  $10^{-15}$ , i.e.,  $j = 15$ . The relationships among relative errors, perturbation bounds, and relative condition numbers are shown in Table 2. Due to the singularity of the matrix  $R$  caused by parameter  $r$ , the accuracy of the perturbation bounds is highly affected by the singularity. When the value of  $m$  increases, the perturbation bound is still tight to the relative error. Also, it can be seen that the number of accurate digits of the perturbation bounds is re-

**Table 2** Relative errors, perturbation bounds and relative condition numbers

$m$	Relative error	$\frac{\xi_*}{\ X\ _F}$	$c_{\text{rel}}(X)$
1	$5.94 \times 10^{-15}$	$1.64 \times 10^{-14}$	$5.84 \times 10^1$
2	$5.28 \times 10^{-14}$	$9.47 \times 10^{-14}$	$5.56 \times 10^2$
3	$4.75 \times 10^{-13}$	$7.85 \times 10^{-13}$	$5.56 \times 10^3$
4	$4.58 \times 10^{-12}$	$7.40 \times 10^{-12}$	$5.56 \times 10^4$
5	$4.85 \times 10^{-11}$	$7.26 \times 10^{-11}$	$5.56 \times 10^5$
6	$4.56 \times 10^{-10}$	$7.22 \times 10^{-10}$	$5.57 \times 10^6$
7	$4.69 \times 10^{-9}$	$8.61 \times 10^{-9}$	$5.57 \times 10^7$

duced proportionally to the increase of the quantities of the relative condition numbers. In other words, if the accurate digits of the perturbation bound are added to the digits in the relative condition numbers, this number is almost equal to 16. (While using IEEE double-precision, the machine precision is around  $2.2 \times 10^{-16}$ .) This implies that the derived perturbation bound of (55) is fairly sharp.

## 7 Conclusion

While doing numerical computation, it is important in practice to have an accurate method for estimating the relative error and the condition number of the given problems. In this paper, we focus on providing a tight perturbation bound of the stabilizing solution to SARE (1a)-(1b) under small changes in the coefficient matrices. Also, some sufficient conditions are presented for the existence of the stabilizing solution to the perturbed SARE. The corresponding condition number of the stabilizing solution is provided in this work. We highlight and compare the practical performance of the derived perturbation bound and condition number through a numerical example. Numerical results show that our perturbation bound is very sensitive to the condition number of the stabilizing solution. As a consequence, they provide good measurement tools for the sensitivity analysis of SARE (1a)-(1b).

## Appendix

We provide here a proof of the condition given by (63).

**Lemma A.1** *Let  $\mathcal{P}, \mathcal{Q}, \mathcal{M}, \mathcal{N}, \mathcal{H}, \mathcal{T}, \mathcal{L}_c^{-1}$  be the operators defined by (58a)-(58b), (24a)-(24d) and (4), and let  $\Omega_\delta, \delta_p$  be defined by (57). Then the following equality holds:*

$$\begin{aligned} & \limsup_{\delta \rightarrow 0^+} \sup_{\Omega_\delta} \frac{\|\mathcal{P}\Delta A + \mathcal{Q}\Delta B + \mathcal{M}\Delta C + \mathcal{N}\Delta D + \mathcal{H}\Delta S + \mathcal{T}\Delta R + \mathcal{L}_c^{-1}\Delta H\|_F}{\delta} \\ &= \max_{\delta_p > 0} \frac{\|\mathcal{P}\Delta A + \mathcal{Q}\Delta B + \mathcal{M}\Delta C + \mathcal{N}\Delta D + \mathcal{H}\Delta S + \mathcal{T}\Delta R + \mathcal{L}_c^{-1}\Delta H\|_F}{\delta_p}. \end{aligned} \quad (65)$$

*Proof* For any  $\delta > 0, 0 < \delta_p \leq \delta$ , we see that

$$\begin{aligned} & \frac{\|\mathcal{P}\Delta A + \mathcal{Q}\Delta B + \mathcal{M}\Delta C + \mathcal{N}\Delta D + \mathcal{H}\Delta S + \mathcal{T}\Delta R + \mathcal{L}_c^{-1}\Delta H\|_F}{\delta} \\ &= \frac{\delta_p}{\delta} \frac{\|\mathcal{P}\Delta A + \mathcal{Q}\Delta B + \mathcal{M}\Delta C + \mathcal{N}\Delta D + \mathcal{H}\Delta S + \mathcal{T}\Delta R + \mathcal{L}_c^{-1}\Delta H\|_F}{\delta_p} \\ &\leq \max_{\delta_{\tilde{p}} > 0} \frac{\|\mathcal{P}\tilde{\Delta A} + \mathcal{Q}\tilde{\Delta B} + \mathcal{M}\tilde{\Delta C} + \mathcal{N}\tilde{\Delta D} + \mathcal{H}\tilde{\Delta S} + \mathcal{T}\tilde{\Delta R} + \mathcal{L}_c^{-1}\tilde{\Delta H}\|_F}{\delta_{\tilde{p}}}, \end{aligned}$$

where  $\delta_{\tilde{p}} = \|(\frac{\tilde{\Delta A}}{\kappa_A}, \frac{\tilde{\Delta B}}{\kappa_B}, \frac{\tilde{\Delta C}}{\kappa_C}, \frac{\tilde{\Delta D}}{\kappa_D}, \frac{\tilde{\Delta S}}{\kappa_S}, \frac{\tilde{\Delta R}}{\kappa_R}, \frac{\tilde{\Delta H}}{\kappa_H})\|_F$ . It follows that

$$\begin{aligned} & \limsup_{\delta \rightarrow 0^+} \sup_{\Omega_\delta} \frac{\|\mathcal{P}\Delta A + \mathcal{Q}\Delta B + \mathcal{M}\Delta C + \mathcal{N}\Delta D + \mathcal{H}\Delta S + \mathcal{T}\Delta R + \mathcal{L}_c^{-1}\Delta H\|_F}{\delta} \\ & \leq \max_{\delta_p > 0} \frac{\|\mathcal{P}\Delta A + \mathcal{Q}\Delta B + \mathcal{M}\Delta C + \mathcal{N}\Delta D + \mathcal{H}\Delta S + \mathcal{T}\Delta R + \mathcal{L}_c^{-1}\Delta H\|_F}{\delta_p}. \end{aligned} \quad (66)$$

On the other hand, for any fixed  $\delta > 0$ , choose any perturbation matrices

$$(\Delta A_1, \Delta B_1, \Delta C_1, \Delta D_1, \Delta S_1, \Delta R_1, \Delta H_1) \in \Omega_\delta$$

and therefore

$$\left\| \left( \frac{\Delta A_1}{\kappa_A}, \frac{\Delta B_1}{\kappa_B}, \frac{\Delta C_1}{\kappa_C}, \frac{\Delta D_1}{\kappa_D}, \frac{\Delta S_1}{\kappa_S}, \frac{\Delta R_1}{\kappa_R}, \frac{\Delta H_1}{\kappa_H} \right) \right\|_F = \delta_{p_1} \leq \delta.$$

It is true that  $(\frac{\delta}{\delta_{p_1}}\Delta A_1, \frac{\delta}{\delta_{p_1}}\Delta B_1, \frac{\delta}{\delta_{p_1}}\Delta C_1, \frac{\delta}{\delta_{p_1}}\Delta D_1, \frac{\delta}{\delta_{p_1}}\Delta S_1, \frac{\delta}{\delta_{p_1}}\Delta R_1, \frac{\delta}{\delta_{p_1}}\Delta H_1) \in \Omega_\delta$  and this gives the fact that

$$\begin{aligned} & \limsup_{\delta \rightarrow 0^+} \sup_{\Omega_\delta} \frac{\|\mathcal{P}\Delta A + \mathcal{Q}\Delta B + \mathcal{M}\Delta C + \mathcal{N}\Delta D + \mathcal{H}\Delta S + \mathcal{T}\Delta R + \mathcal{L}_c^{-1}\Delta H\|_F}{\delta} \\ & \geq \frac{\|\mathcal{P}\frac{\delta}{\delta_{p_1}}\Delta A_1 + \mathcal{Q}\frac{\delta}{\delta_{p_1}}\Delta B_1 + \mathcal{M}\frac{\delta}{\delta_{p_1}}\Delta C_1 + \mathcal{N}\frac{\delta}{\delta_{p_1}}\Delta D_1 + \mathcal{H}\frac{\delta}{\delta_{p_1}}\Delta S_1 + \mathcal{T}\frac{\delta}{\delta_{p_1}}\Delta R_1 + \mathcal{L}_c^{-1}\frac{\delta}{\delta_{p_1}}\Delta H_1\|_F}{\delta} \\ & = \frac{\|\mathcal{P}\Delta A_1 + \mathcal{Q}\Delta B_1 + \mathcal{M}\Delta C_1 + \mathcal{N}\Delta D_1 + \mathcal{H}\Delta S_1 + \mathcal{T}\Delta R_1 + \mathcal{L}_c^{-1}\Delta H_1\|_F}{\delta_{p_1}}. \end{aligned}$$

Hence

$$\begin{aligned} & \limsup_{\delta \rightarrow 0^+} \sup_{\Omega_\delta} \frac{\|\mathcal{P}\Delta A + \mathcal{Q}\Delta B + \mathcal{M}\Delta C + \mathcal{N}\Delta D + \mathcal{H}\Delta S + \mathcal{T}\Delta R + \mathcal{L}_c^{-1}\Delta H\|_F}{\delta} \\ & \geq \max_{\delta_p > 0} \frac{\|\mathcal{P}\Delta A + \mathcal{Q}\Delta B + \mathcal{M}\Delta C + \mathcal{N}\Delta D + \mathcal{H}\Delta S + \mathcal{T}\Delta R + \mathcal{L}_c^{-1}\Delta H\|_F}{\delta_p}. \end{aligned} \quad (67)$$

Comparison of (66) and (67) gives (65).  $\square$

#### Competing interests

The authors declare that there is no conflict of interests regarding the publication of this article.

#### Authors' contributions

All authors contributed equally and significantly in writing this paper. All authors read and approved the final manuscript.

#### Author details

<sup>1</sup>Center for General Education, National Formosa University, Huwei 632, Taiwan. <sup>2</sup>Department of Mathematics, National Taiwan Normal University, Taipei 116, Taiwan. <sup>3</sup>Department of Mathematics, National Chung Cheng University, Chia-Yi 621, Taiwan.

#### Acknowledgements

The authors wish to thank the editor and two anonymous referees for many interesting and valuable suggestions on the manuscript. This research work is partially supported by the National Science Council and the National Center for Theoretical Sciences in Taiwan. The first author was supported by the National Science Council of Taiwan under Grant NSC 102-2115-M-150-002. The second author was supported by the National Science Council of Taiwan under Grant NSC 102-2115-M-003-009. The third author was supported by the National Science Council of Taiwan under Grant NSC 101-2115-M-194-007-MY3.

Received: 16 September 2013 Accepted: 14 November 2013 Published: 11 Dec 2013

# References

1. Rami, MA, Zhou, XY: Linear matrix inequalities, Riccati equations, and indefinite stochastic linear quadratic controls. *IEEE Trans. Autom. Control* **45**(6), 1131-1143 (2000). doi:10.1109/9.863597
2. El Bouhtouri, A, Hinrichsen, D, Pritchard, AJ: On the disturbance attenuation problem for a wide class of time invariant linear stochastic systems. *Stoch. Stoch. Rep.* **65**(3-4), 255-297 (1999)
3. Damm, T, Hinrichsen, D: Newton's method for a rational matrix equation occurring in stochastic control. In: *Proceedings of the Eighth Conference of the International Linear Algebra Society (Barcelona, 1999)*, vol. 332/334, pp. 81-109 (2001)
4. Hinrichsen, D, Pritchard, AJ: Stochastic  $H_\infty$ . *SIAM J. Control Optim.* **36**, 1504-1538 (1998)
5. Lancaster, P, Rodman, L: *Algebraic Riccati Equations*. Oxford Science Publications. Clarendon, New York (1995)
6. Mehrmann, VL: *The Autonomous Linear Quadratic Control Problem: Theory and Numerical Solution*. Lecture Notes in Control and Information Sciences, vol. 163. Springer, Berlin (1991). doi:10.1007/BFb0039443
7. Zhou, K, Doyle, JC, Glover, K: *Robust and Optimal Control*. Prentice Hall, Upper Saddle River (1996)
8. Benner, P, Laub, AJ, Mehrmann, V: A collection of benchmark examples for the numerical solution of algebraic Riccati equations I: continuous-time case. Technical Report SPC 95\_22, Fakultät für Mathematik, TU Chemnitz-Zwickau, 09107 Chemnitz, FRG. <http://www.tu-chemnitz.de/sfb393/spc95pr.html> (1995)
9. Benner, P, Laub, AJ, Mehrmann, V: A collection of benchmark examples for the numerical solution of algebraic Riccati equations II: discrete-time case. Technical Report SPC 95\_23, Fakultät für Mathematik, TU Chemnitz-Zwickau, 09107 Chemnitz, FRG. <http://www.tu-chemnitz.de/sfb393/spc95pr.html> (1995)
10. Sima, V: *Algorithms for Linear-Quadratic Optimization*. Monographs and Textbooks in Pure and Applied Mathematics, vol. 200. Dekker, New York (1996)
11. Laub, AJ: A Schur method for solving algebraic Riccati equations. *IEEE Trans. Autom. Control* **24**(6), 913-921 (1979)
12. Ammar, G, Benner, P, Mehrmann, V: A multishift algorithm for the numerical solution of algebraic Riccati equations. *Electron. Trans. Numer. Anal.* **1**, 33-48 (1993)
13. Ammar, G, Mehrmann, V: On Hamiltonian and symplectic Hessenberg forms. *Linear Algebra Appl.* **149**, 55-72 (1991)
14. Benner, P, Mehrmann, V, Xu, H: A new method for computing the stable invariant subspace of a real Hamiltonian matrix. *J. Comput. Appl. Math.* **86**, 17-43 (1997)
15. Benner, P, Mehrmann, V, Xu, H: A numerically stable, structure preserving method for computing the eigenvalues of real Hamiltonian or symplectic pencils. *Numer. Math.* **78**(3), 329-358 (1998)
16. Bunse-Gerstner, A, Byers, R, Mehrmann, V: A chart of numerical methods for structured eigenvalue problems. *SIAM J. Matrix Anal. Appl.* **13**, 419-453 (1992)
17. Bunse-Gerstner, A, Mehrmann, V: A symplectic QR like algorithm for the solution of the real algebraic Riccati equation. *IEEE Trans. Autom. Control* **31**(12), 1104-1113 (1986)
18. Mehrmann, V: A step toward a unified treatment of continuous and discrete time control problems. *Linear Algebra Appl.* **241-243**, 749-779 (1996)
19. Byers, R: Solving the algebraic Riccati equation with the matrix sign function. *Linear Algebra Appl.* **85**, 267-279 (1987)
20. Benner, P: Contributions to the numerical solutions of algebraic Riccati equations and related eigenvalue problems. PhD thesis, Fakultät für Mathematik, TU Chemnitz-Zwickau, Chemnitz, Germany (1997)
21. Chu, EK-W, Fan, H-Y, Lin, W-W: A structure-preserving doubling algorithm for continuous-time algebraic Riccati equations. *Linear Algebra Appl.* **396**, 55-80 (2005)
22. Chu, EK-W, Fan, H-Y, Lin, W-W, Wang, C-S: Structure-preserving algorithms for periodic discrete-time algebraic Riccati equations. *Int. J. Control* **77**(8), 767-788 (2004)
23. Chiang, C-Y, Fan, H-Y: Residual bounds of the stochastic algebraic Riccati equation. *Appl. Numer. Math.* **63**, 78-87 (2013)
24. Konstantinov, M, Gu, D-W, Mehrmann, V, Petkov, P: *Perturbation Theory for Matrix Equations*. Studies in Computational Mathematics, vol. 9. North-Holland, Amsterdam (2003)
25. Sun, J-G: Perturbation theory for algebraic Riccati equations. *SIAM J. Matrix Anal. Appl.* **19**(1), 39-65 (1998). doi:10.1137/S0895479895291303
26. Sun, J-g: Sensitivity analysis of the discrete-time algebraic Riccati equation. In: *Proceedings of the Sixth Conference of the International Linear Algebra Society (Chemnitz, 1996)*, vol. 275/276, pp. 595-615 (1998). doi:10.1016/S0024-3795(97)10017-9
27. Sun, J-g: Residual bounds of approximate solutions of the algebraic Riccati equation. *Numer. Math.* **76**(2), 249-263 (1997). doi:10.1007/s002110050262
28. Sun, J-g: Residual bounds of approximate solutions of the discrete-time algebraic Riccati equation. *Numer. Math.* **78**(3), 463-478 (1998). doi:10.1007/s002110050321
29. Riedel, KS: A Sherman-Morrison-Woodbury identity for rank augmenting matrices with application to centering. *SIAM J. Matrix Anal. Appl.* **13**(2), 659-662 (1992). doi:10.1137/0613040
30. Stewart, GW, Sun, JG: *Matrix Perturbation Theory*. Computer Science and Scientific Computing. Academic Press, Boston (1990)
31. Ortega, JM, Rheinboldt, WC: *Iterative Solution of Nonlinear Equations in Several Variables*. Classics in Applied Mathematics, vol. 30. SIAM, Philadelphia (2000). Reprint of the 1970 original
32. Rice, JR: A theory of condition. *SIAM J. Numer. Anal.* **3**, 287-310 (1966)
33. Sun, J-g: Condition numbers of algebraic Riccati equations in the Frobenius norm. *Linear Algebra Appl.* **350**, 237-261 (2002). doi:10.1016/S0024-3795(02)00294-X
34. Xu, S: *Matrix Computation in Control Theory*. Higher Education Press, Beijing (2010) (In Chinese)

10.1186/1029-242X-2013-580

**Cite this article as:** Chiang et al.: Perturbation analysis of the stochastic algebraic Riccati equation. *Journal of Inequalities and Applications* 2013, **2013**:580